

Predicting time-to-event outcomes based on high-dimensional multivariate longitudinal information

EFSPI: Brussels, November 7, 2013

Geert Verbeke

geert.verbeke@med.kuleuven.be

http://perswww.kuleuven.be/geert_verbeke



KU LEUVEN

Interuniversity Institute for Biostatistics
and statistical Bioinformatics

Motivating example: Renal graft failure

- Patients with kidney transplant between 1983 and 2000 at U.H.Leuven
- Clinical interest:

Continuous prediction of long-term success of graft (> 10 years)

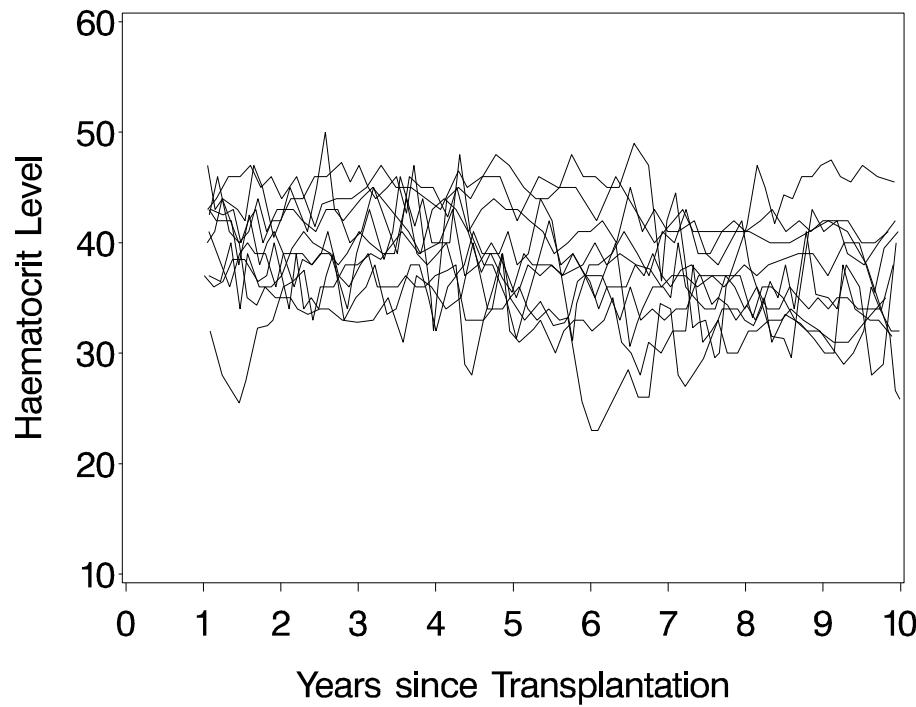
- Conditional on:
 - ▷ not losing graft during first year
 - ▷ not dying in the first 10 years for reasons not related to transplantation.

Motivating example (cont'd)

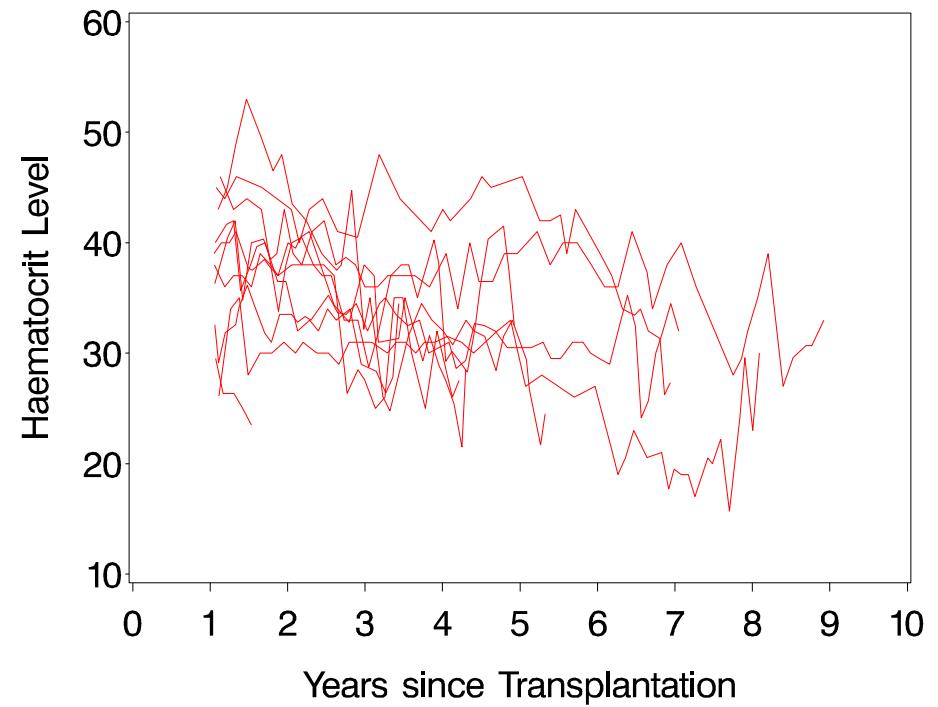
- Information: 949 patients, with 1-78 visits per patient
 - ▷ 341 patients with functioning graft after 10 years
 - ▷ 91 patients with a graft failure before 10 years
 - ▷ 517 patients with functioning graft, but FU < 10 yrs
- Prediction based on longitudinal measurements of:
 - ▷ Haematocrit Level
 - ▷ Filtration Rate
 - ▷ Proteinuria
 - ▷ Blood Pressure

Haematocrit level

Non-failures

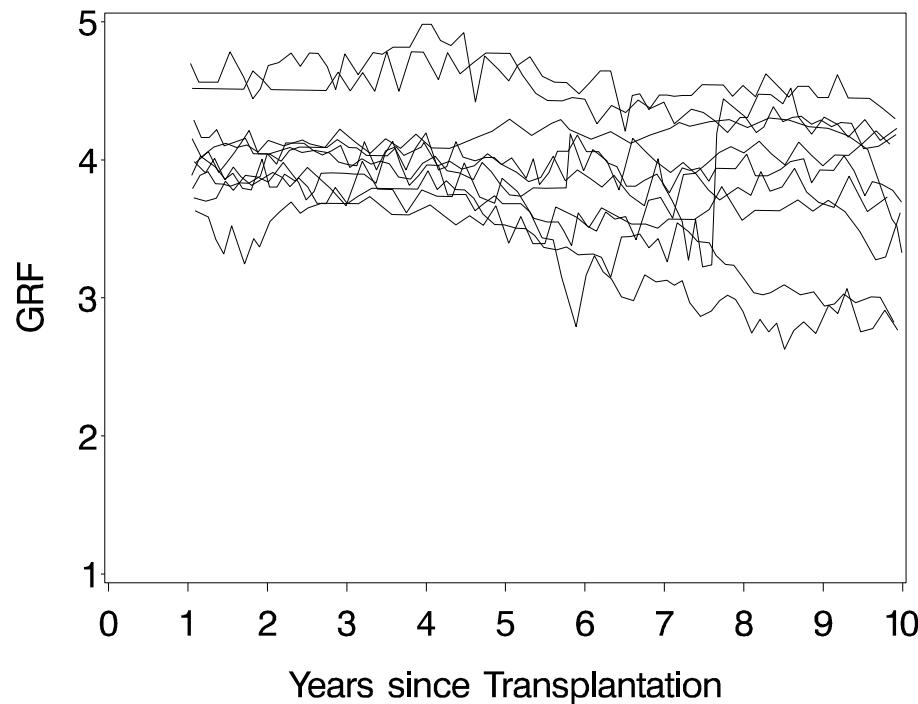


Failures

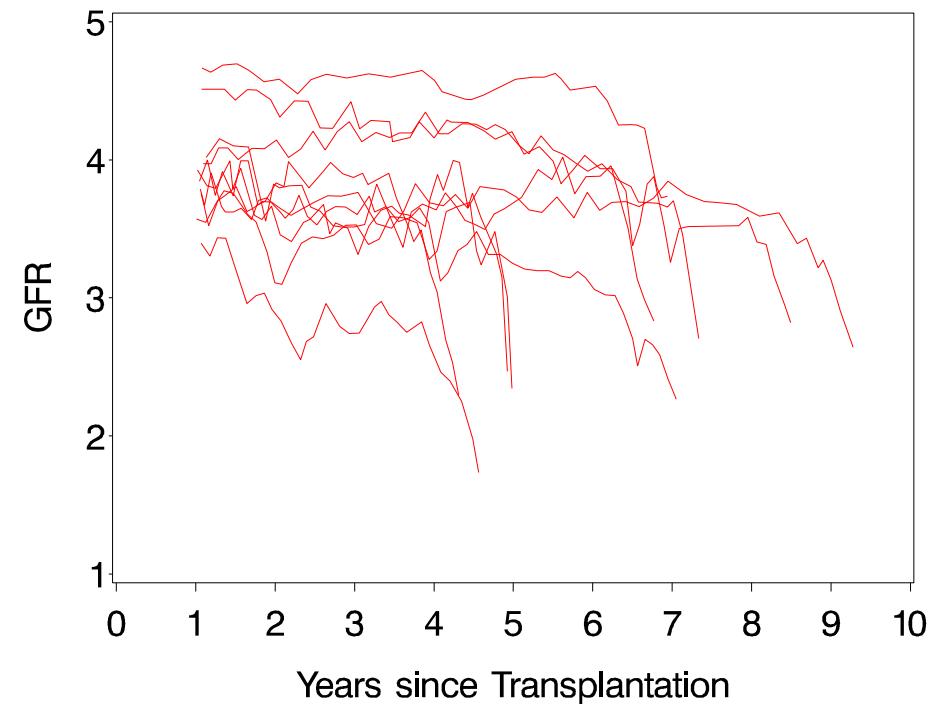


Glomerular filtration rate (GFR)

Non-failures

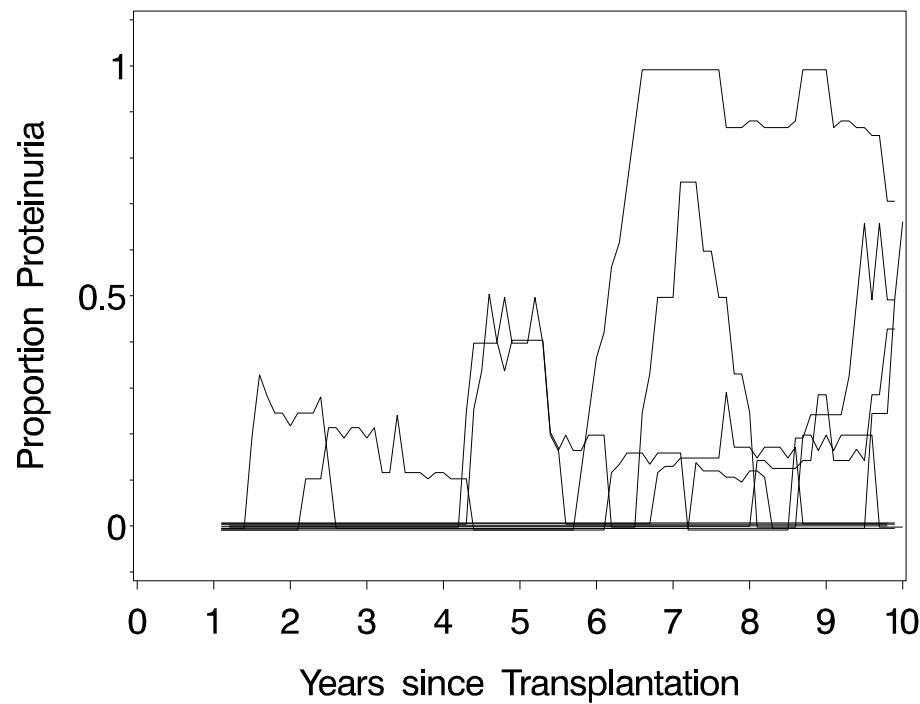


Failures

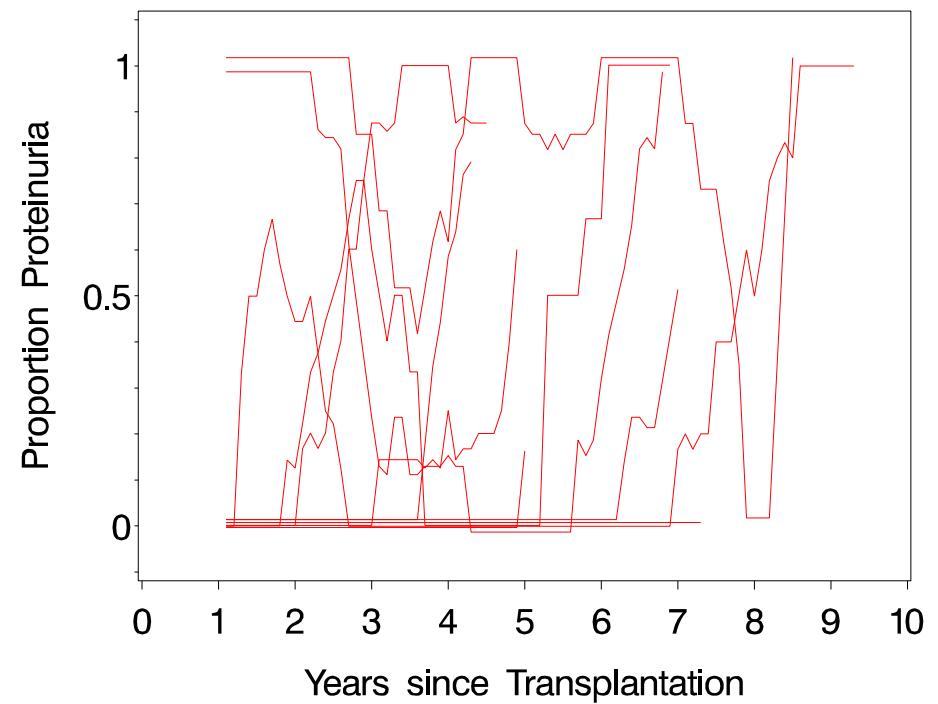


Presence of proteinuria

Non-failures

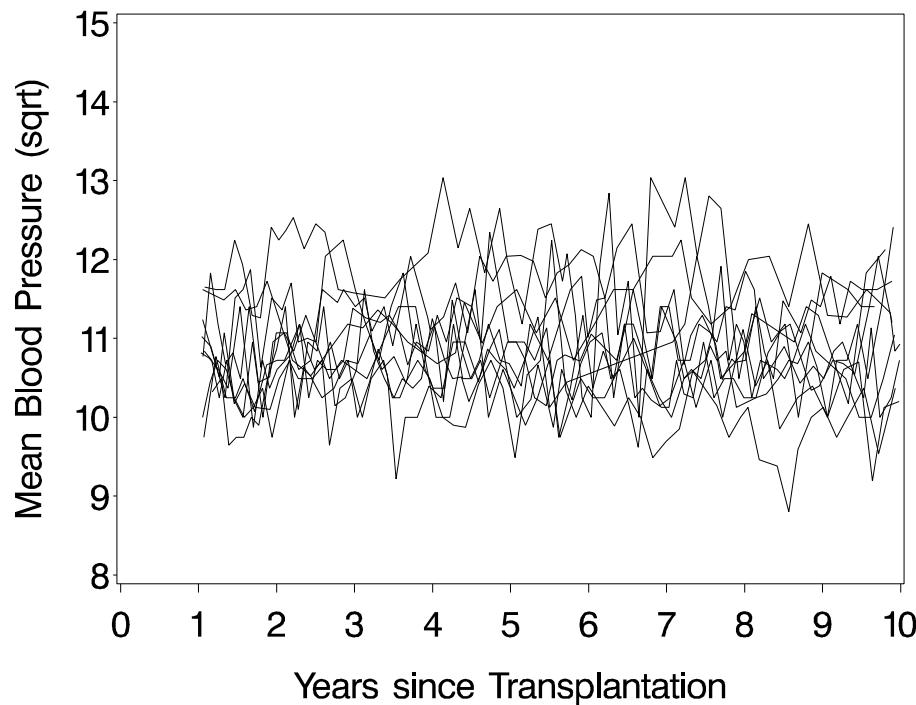


Failures

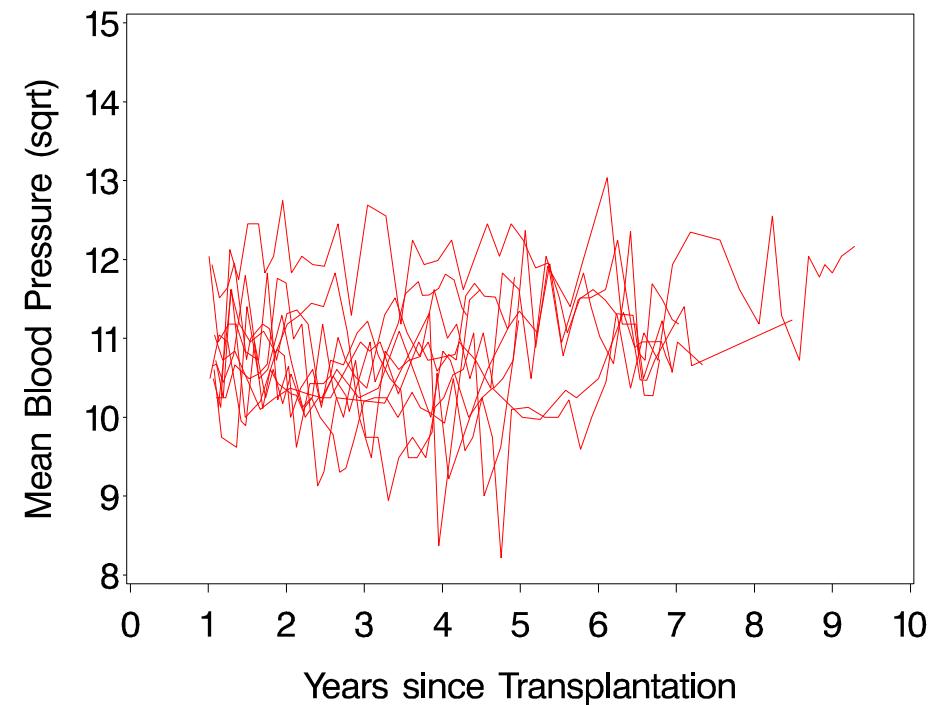


Mean blood pressure

Non-failures



Failures



Aim of the analysis

$$H_i(t) = P_i(t \leq F_i \leq 120 \mid \mathbf{y}_i^{\leq t}), \quad \forall t$$

Specification of conditional distribution for $(F_i \mid \mathbf{Y}_i^{\leq t})$ problematic due to:

- ▷ Unbalanced nature of the longitudinal data
- ▷ The different outcome types in \mathbf{Y}_i
- ▷ The running time t

A pattern-mixture approach

$$f(F_i, \mathbf{Y}_i) = f(\mathbf{Y}_i|F_i)f(F_i)$$



$$f(F_i, \mathbf{Y}_i^{\leq t}) = f(\mathbf{Y}_i^{\leq t}|F_i)f(F_i)$$



+ Bayes rule

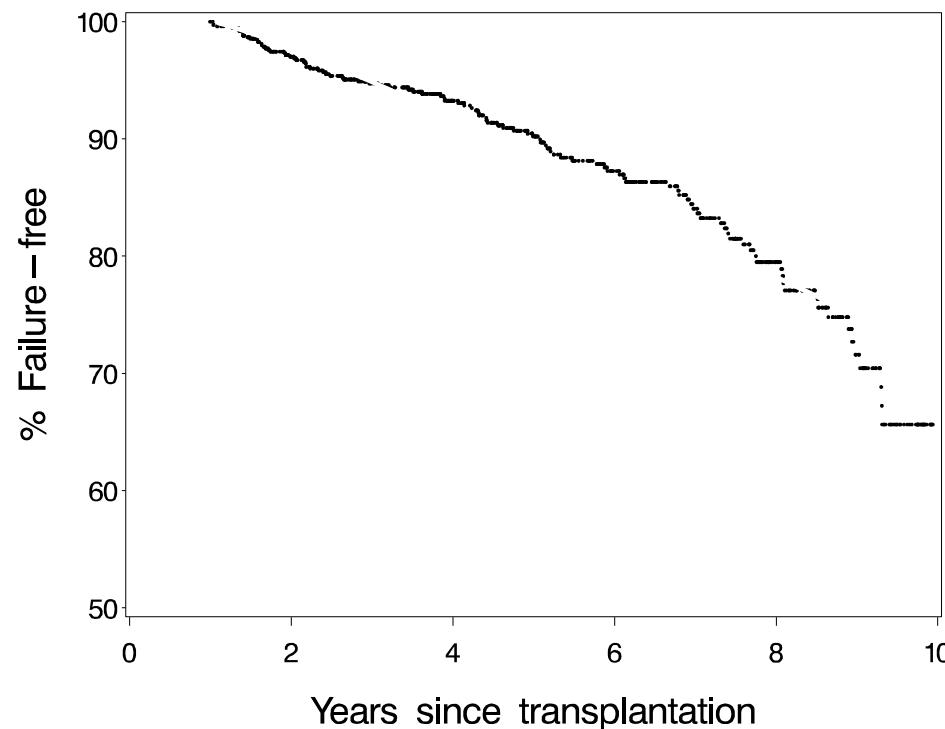
$$H_i(t) = P_i(t \leq F_i \leq 120 \mid \mathbf{y}_i^{\leq t})$$

$$= \frac{f_i(\mathbf{y}_i^{\leq t} \mid t \leq F_i \leq 120)P(F_i \leq 120 \mid F_i \geq t)}{f_i(\mathbf{y}_i^{\leq t} \mid t \leq F_i \leq 120)P(F_i \leq 120 \mid F_i \geq t) + f_i(\mathbf{y}_i^{\leq t} \mid F_i > 120)P(F_i > 120 \mid F_i \geq t)}$$

Prior probabilities

$$H_i(t) = \frac{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t)}{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t) + f_i(\mathbf{y}_i^{\leq t} | F_i > 120) P(F_i > 120 | F_i \geq t)}$$

Time-specific prior probabilities from Kaplan-Meier estimate:



Models for longitudinal outcomes

$$H_i(t) = \frac{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t)}{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t) + f_i(\mathbf{y}_i^{\leq t} | F_i > 120) P(F_i > 120 | F_i \geq t)}$$

- Mixed model for each outcome separately:

$$\mathbf{Y}_{1i} | \mathbf{b}_{1i} \sim G_{1i}(\boldsymbol{\psi}_1, \mathbf{b}_{1i}), \quad \dots, \quad \mathbf{Y}_{4i} | \mathbf{b}_{4i} \sim G_{4i}(\boldsymbol{\psi}_4, \mathbf{b}_{4i})$$

- Linear, generalized linear, and non-linear mixed models possible
- Joint model through joint distribution for all random effects:
$$(\mathbf{b}'_{1i}, \dots, \mathbf{b}'_{4i})' \sim N(\mathbf{0}, D)$$
- Advantage: Model building for each outcome separately

Description of non-failures

$$H_i(t) = \frac{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t)}{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t) + f_i(\mathbf{y}_i^{\leq t} | F_i > 120) P(F_i > 120 | F_i \geq t)}$$

- Outcomes measured during first 10 yrs.
- If failure, then only after 10 yrs.
- Assumption: Models do not depend on F_i

Mixed models for non-failures

- Haematocrit:

$$Y_{1i}(t) = \beta_{01} + b_{01i} + (\beta_{11} + b_{11i})t + \varepsilon_{1i}(t)$$

- GFR:

$$Y_{2i}(t) = \beta_{02} + b_{02i} + (\beta_{12} + b_{12i})t + \varepsilon_{2i}(t)$$

- Proteinuria:

$$\text{logit}\{P(Y_{3i}(t))\} = \beta_{03} + b_{03i} + \beta_{13}t$$

- Mean Blood Pressure:

$$Y_{4i}(t) = \beta_{04} + b_{04i} + (\beta_{14} + b_{14i})t + \varepsilon_{4i}(t)$$

Description of failures

$$H_i(t) = \frac{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t)}{f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) P(F_i \leq 120 | F_i \geq t) + f_i(\mathbf{y}_i^{\leq t} | F_i > 120) P(F_i > 120 | F_i \geq t)}$$

- Outcomes measured until moment of failure
- Implication: Models depend on F_i :

$$\begin{aligned} & f_i(\mathbf{y}_i^{\leq t} | t \leq F_i \leq 120) \\ &= \sum_{k=t}^{119} f_i(\mathbf{y}_i^{\leq t} | k \leq F_i \leq k+1) P(k \leq F_i \leq k+1) / P(t \leq F_i \leq 120) \\ &\approx \sum_{k=t}^{119} f_i(\mathbf{y}_i^{\leq t} | F_i = k + \frac{1}{2}) P(k \leq F_i \leq k+1) / P(t \leq F_i \leq 120) \end{aligned}$$

Mixed models for failures

- Haematocrit:

$$Y_{1i}(t) = \beta_{01} + \gamma_{01}F_i + b_{01i} + (\beta_{11} + \gamma_{11}F_i + b_{11i})t + \varepsilon_{1i}(t)$$

- GFR:

$$Y_{2i}(t) = \begin{cases} \phi_0 + b_{02i} + \beta_{12}[t - \phi_2] + \varepsilon_{2i}(t) & \text{if } t \leq \phi_2 \\ \phi_0 + b_{02i} + \beta_{32}[t - \phi_2] + \varepsilon_{2i}(t) & \text{if } t > \phi_2 \end{cases}$$

with $\phi_0 = \beta_{02} + \gamma_{02}F_i$ and $\phi_2 = \beta_{22} + \gamma_{22}F_i$

- Proteinuria:

$$\text{logit}\{P(Y_{3i}(t))\} = \beta_{03} + \gamma_{03}F_i + b_{03i} + (\beta_{13} + \gamma_{13}F_i)t$$

- Mean Blood Pressure:

$$Y_{4i}(t) = \beta_{04} + \gamma_{04}F_i + b_{04i} + (\beta_{14} + \gamma_{14}F_i + b_{14i})t + \varepsilon_{4i}(t)$$

Mixed models: Summary

	Non-Failures	Failures
Haematocrit:	LMM (2)	LMM (2)
GFR:	LMM (2)	NLMM (1)
Proteinuria:	GLMM (1)	GLMM (1)
Mean Blood Pressure:	LMM (2)	LMM (2)

- ➡ 2 mixed models with many random effects (7 & 6)
- ➡ computational difficulties

Joint mixed model: Pairwise approach

- Fit all 6 bivariate models using (RE)ML:

$$(\mathbf{Y}_1, \mathbf{Y}_2), (\mathbf{Y}_1, \mathbf{Y}_3), (\mathbf{Y}_1, \mathbf{Y}_4), (\mathbf{Y}_2, \mathbf{Y}_3), (\mathbf{Y}_2, \mathbf{Y}_4), (\mathbf{Y}_3, \mathbf{Y}_4)$$

- Equivalent to maximizing pseudo likelihood:

$$p\ell(\boldsymbol{\Theta}) = \ell(\boldsymbol{\Theta}_{1,2} | \mathbf{Y}_1, \mathbf{Y}_2) + \ell(\boldsymbol{\Theta}_{1,3} | \mathbf{Y}_1, \mathbf{Y}_3) + \dots + \ell(\boldsymbol{\Theta}_{3,4} | \mathbf{Y}_3, \mathbf{Y}_4)$$

- Asymptotic properties (from pseudo likelihood theory):

$$\sqrt{N}(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}) \sim MVN(\mathbf{0}, J^{-1}KJ^{-1})$$

J and K consist of first and second-order derivatives of $p\ell$.

- Multiple estimates for same parameters are averaged

Fieuws & Verbeke, Biometrics (2006)

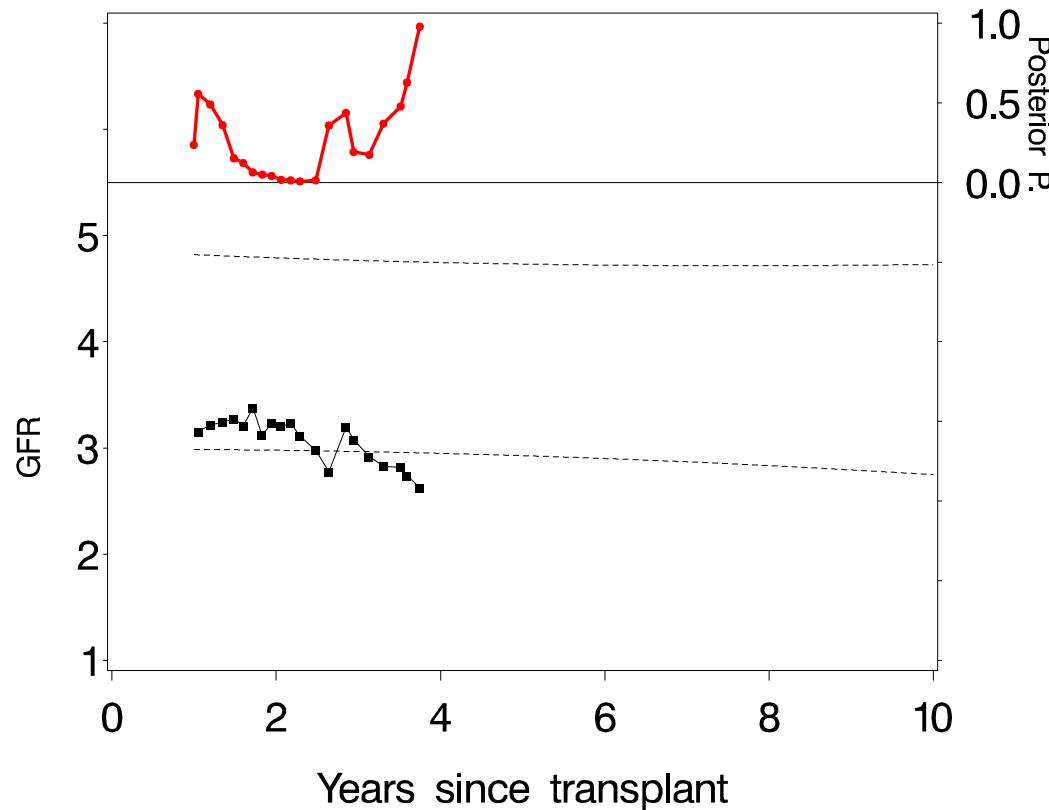
Association between markers ?

L.R. test for no association in each pairwise model:

Markers	Non-failures			Failures				
	Δ	dev.	#	p	Δ	dev.	#	p
1,2:	82.8	4	< 0.0001	18.9	2	< 0.0001		
1,3:	18	2	0.0001	7.2	2	0.027		
1,4:	6.4	4	0.17	2.6	4	0.63		
2,3:	22.4	2	< 0.0001	0.3	1	0.58		
2,4:	8.1	4	0.09	0.1	2	0.95		
3,4:	7.4	2	0.025	6.1	2	0.047		

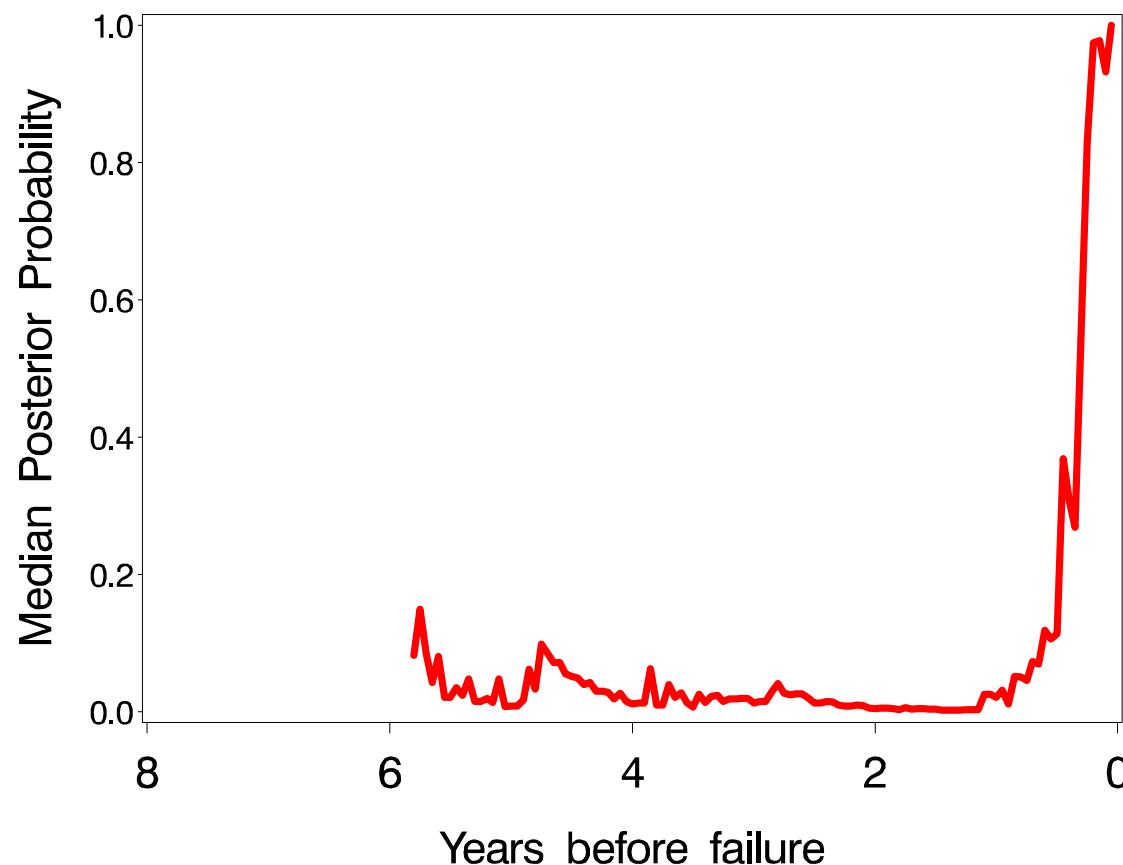
Posterior probabilities based on GFR only

- Training and validation dataset (50% of patients)
- Example for 1 failure from validation dataset:



Median posterior probability, GFR only

- All failures in validation dataset:



Discriminant analysis using 4 markers

Strategies:

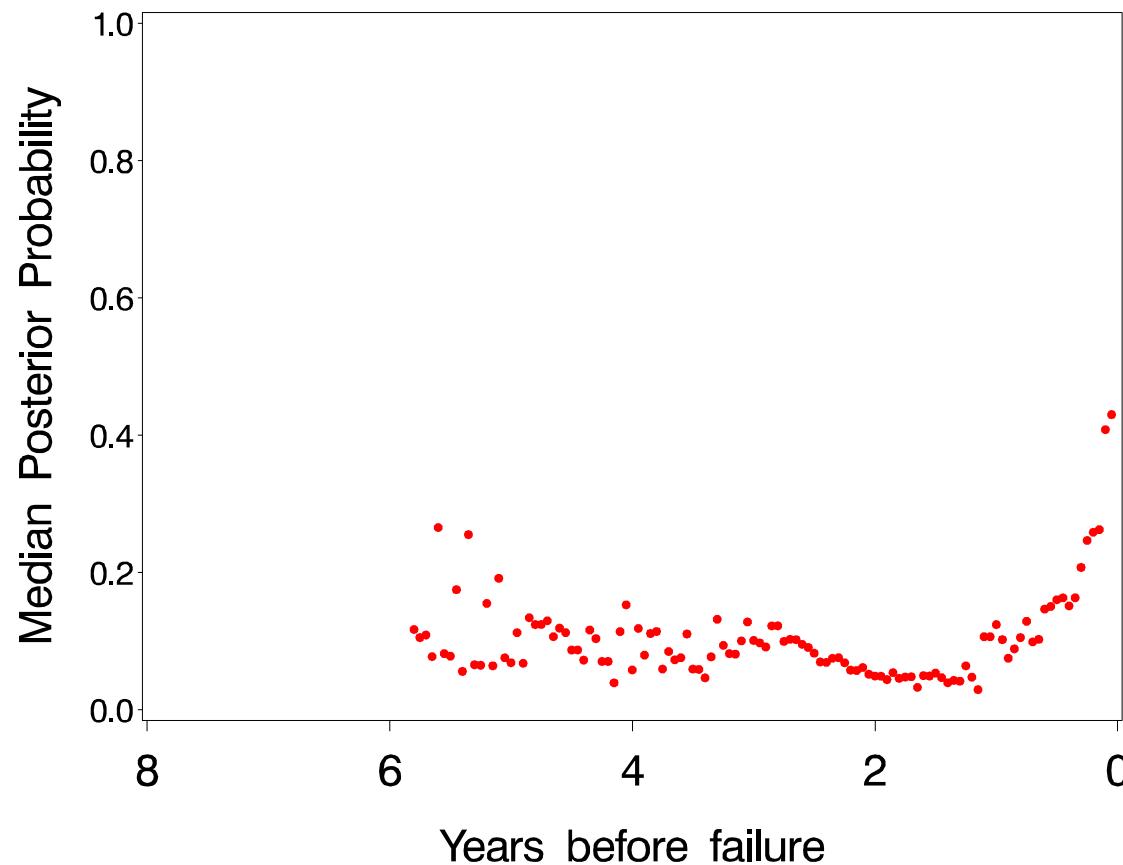
- Decision based on highest posterior probability
- Joint model assuming uncorrelated markers
- Joint model allowing markers to be correlated

Median posterior probability: Failures

- 46 patients in validation set who fail
- Median posterior probabilities to fail within the remaining period
- As a function of time: years before failure

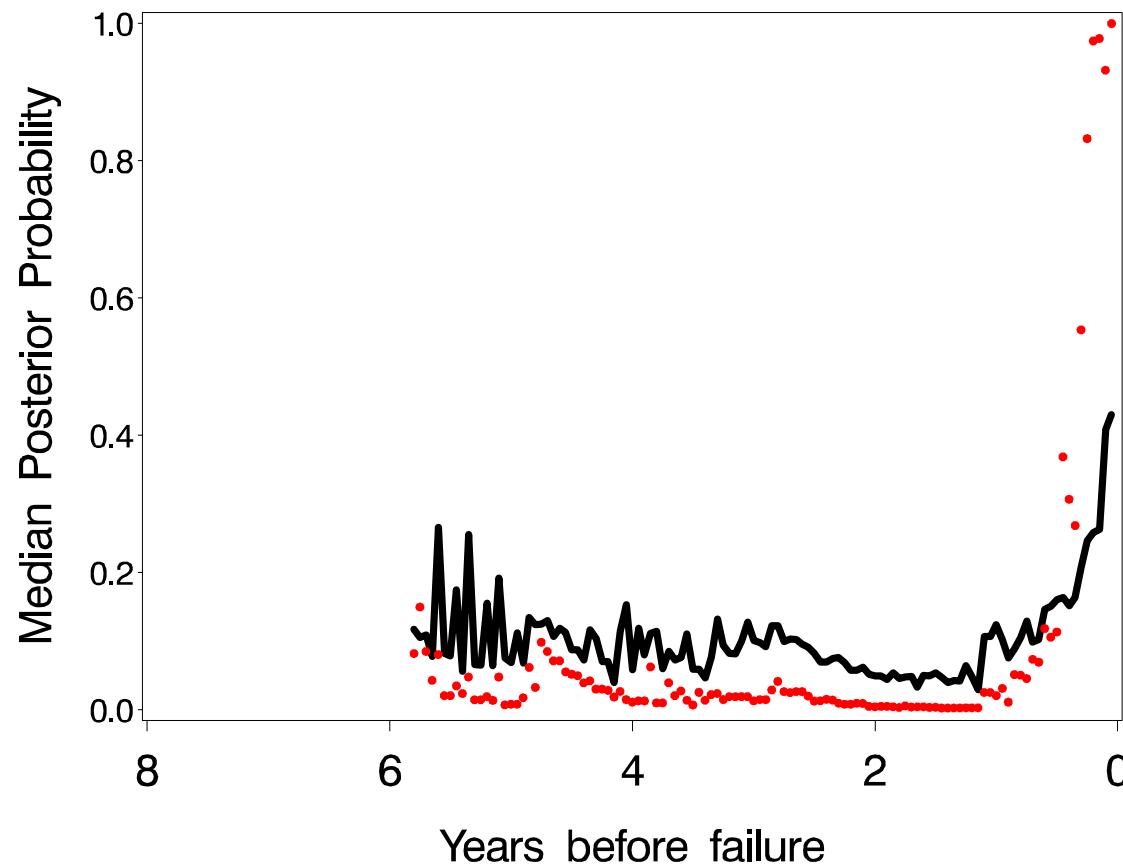
Median posterior probability: Failures

- Only Haematocrit:



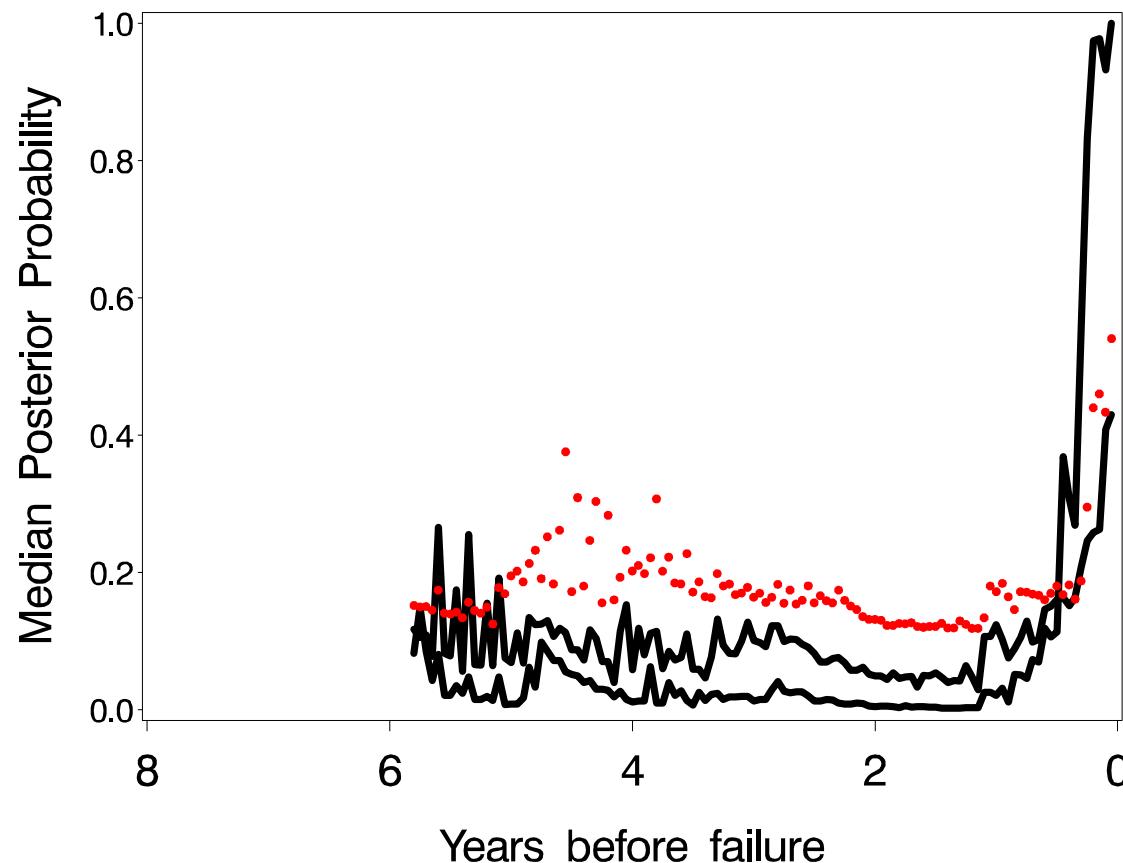
Median posterior probability: Failures

- Only GFR:



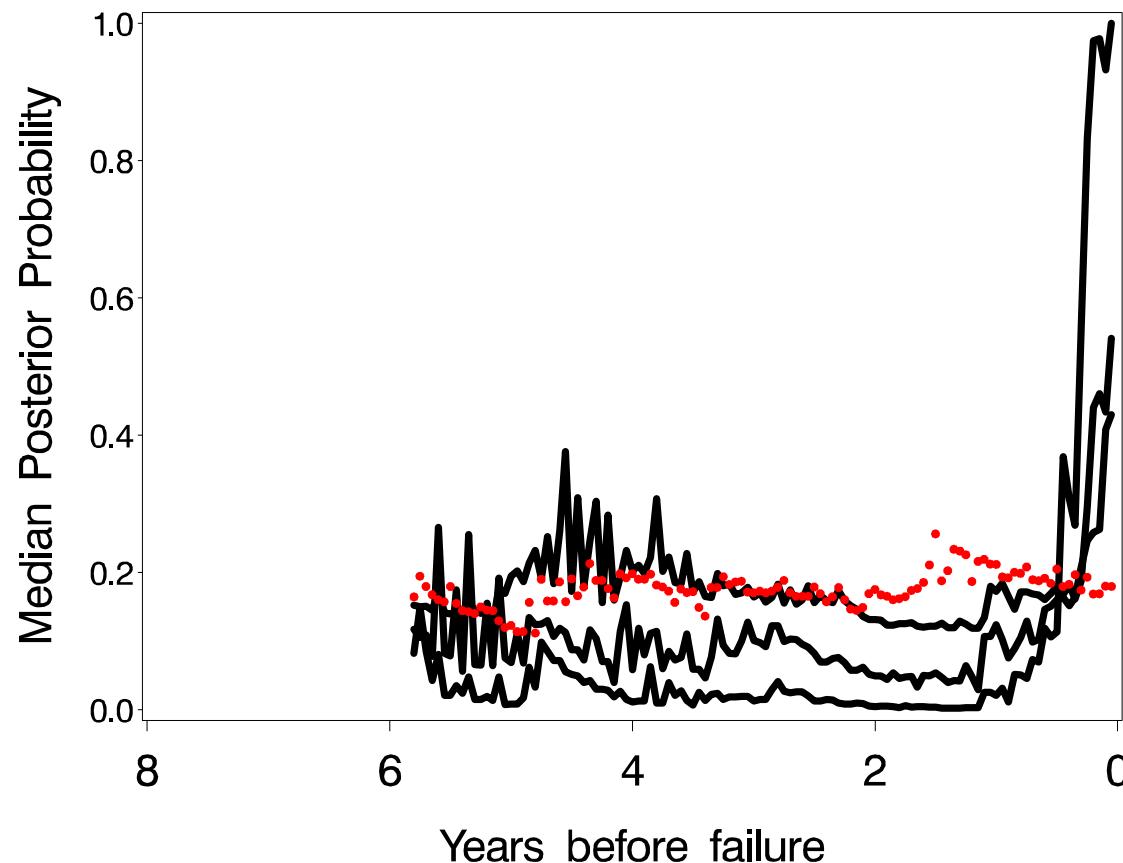
Median posterior probability: Failures

- Only Proteinuria:



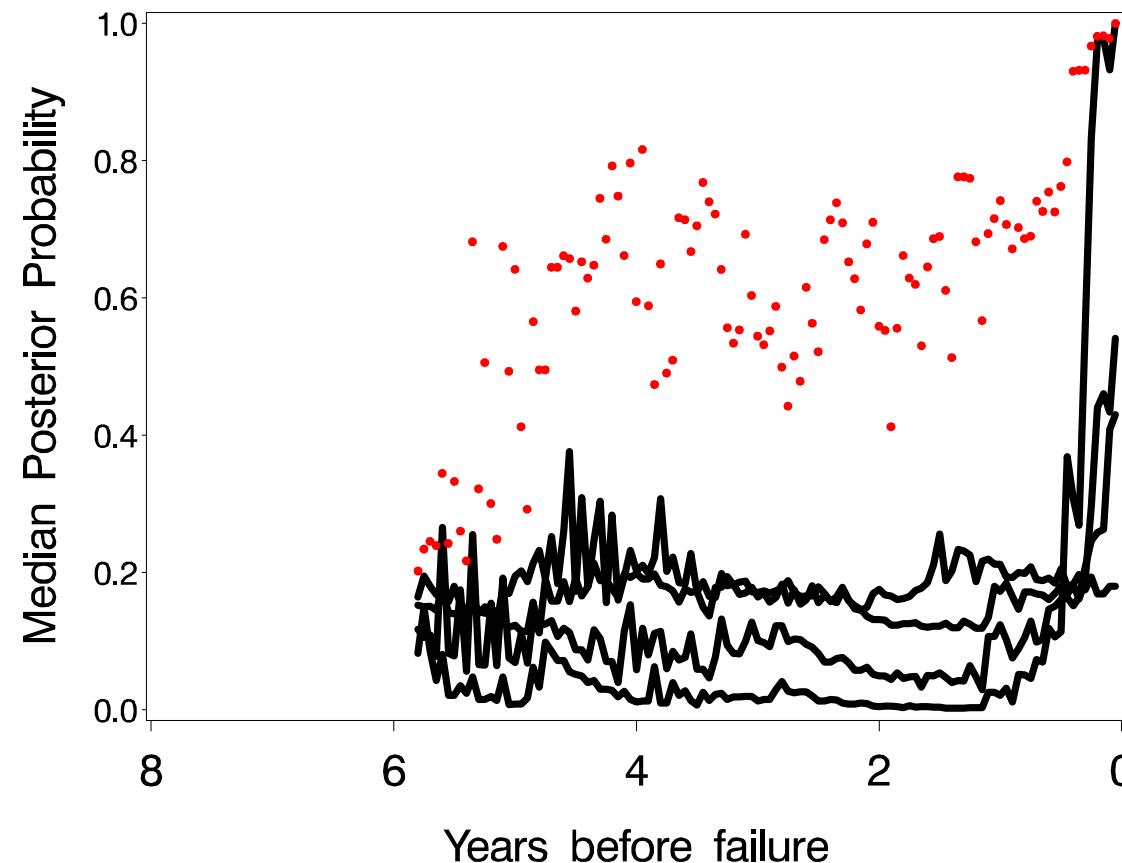
Median posterior probability: Failures

- Only Mean Blood Pressure:



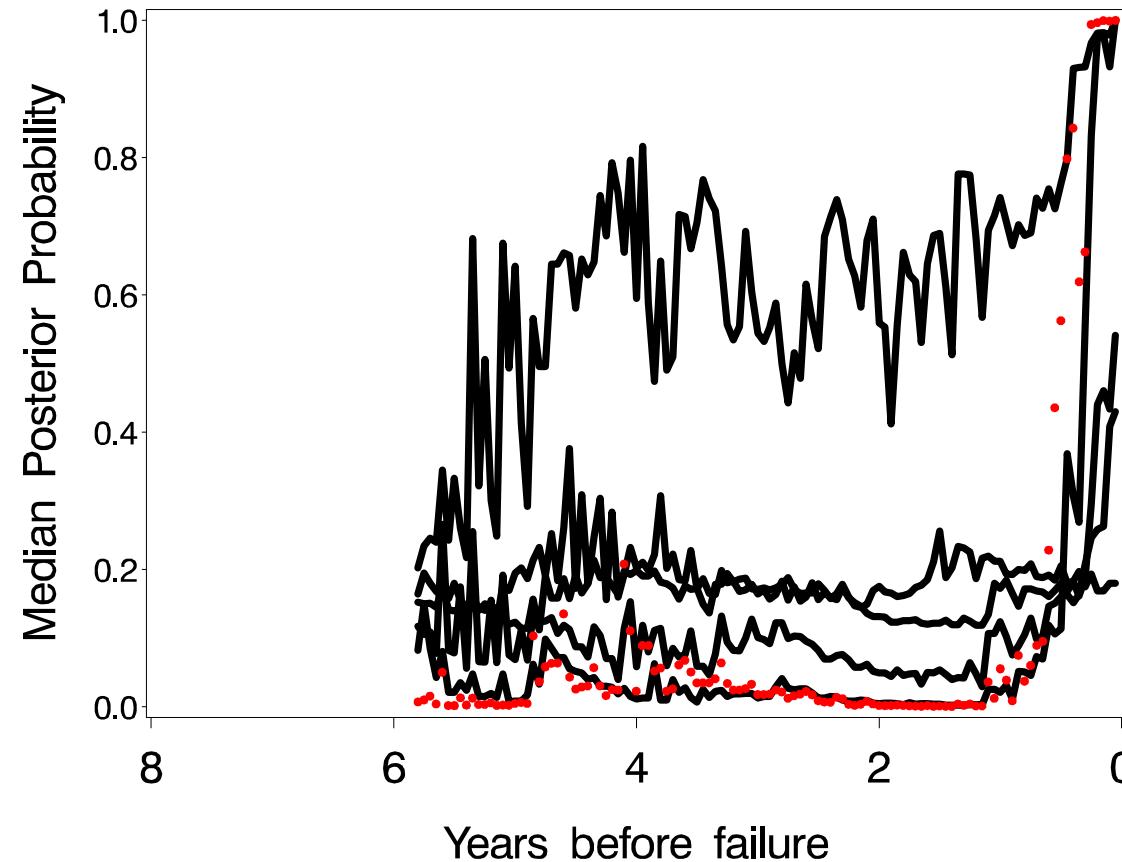
Median posterior probability: Failures

- Highest posterior probability over all four markers:



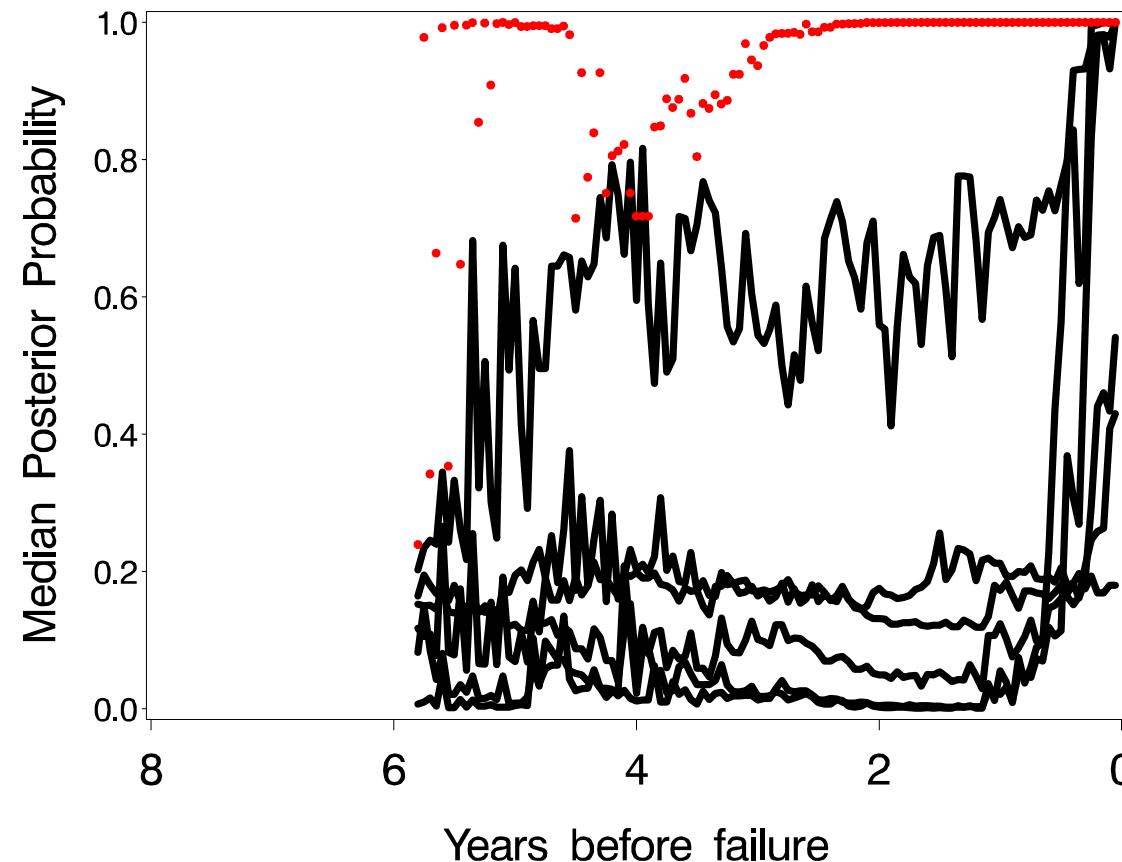
Median posterior probability: Failures

- All four markers, using joint model with uncorrelated markers:



Median posterior probability: Failures

- All four markers, using joint model with correlated markers:

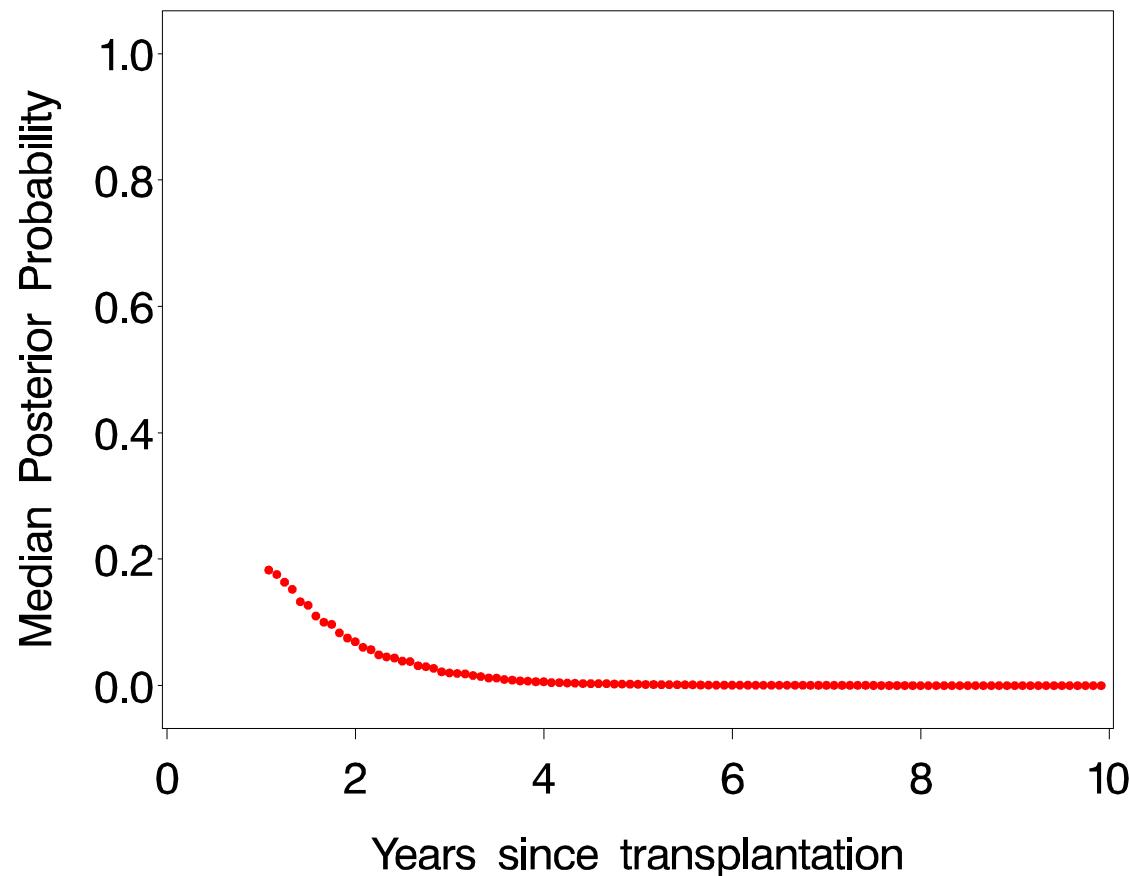


Median posterior probability: Non-failures

- 171 patients in validation set who do not fail
- Median posterior probabilities to fail within the remaining period
- As a function of time: years since transplantation

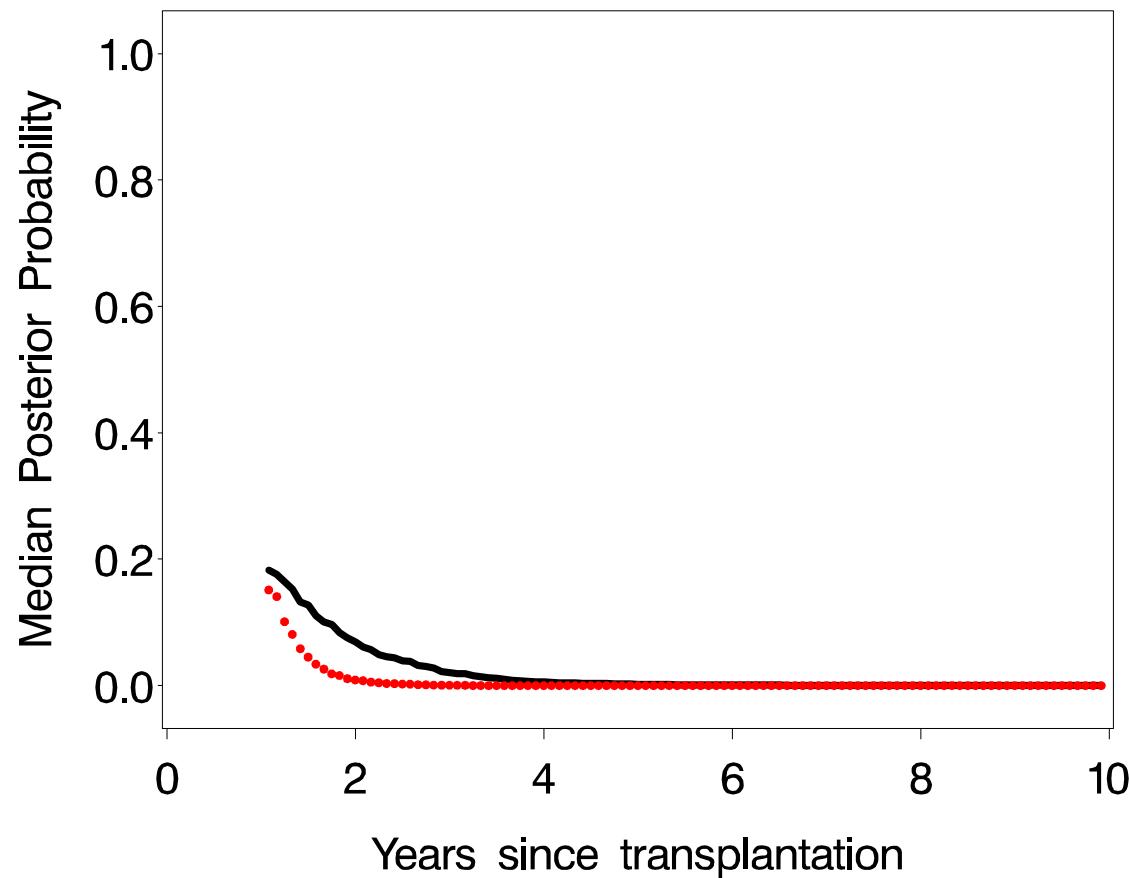
Median posterior probability: Non-failures

- Only Haematocrit:



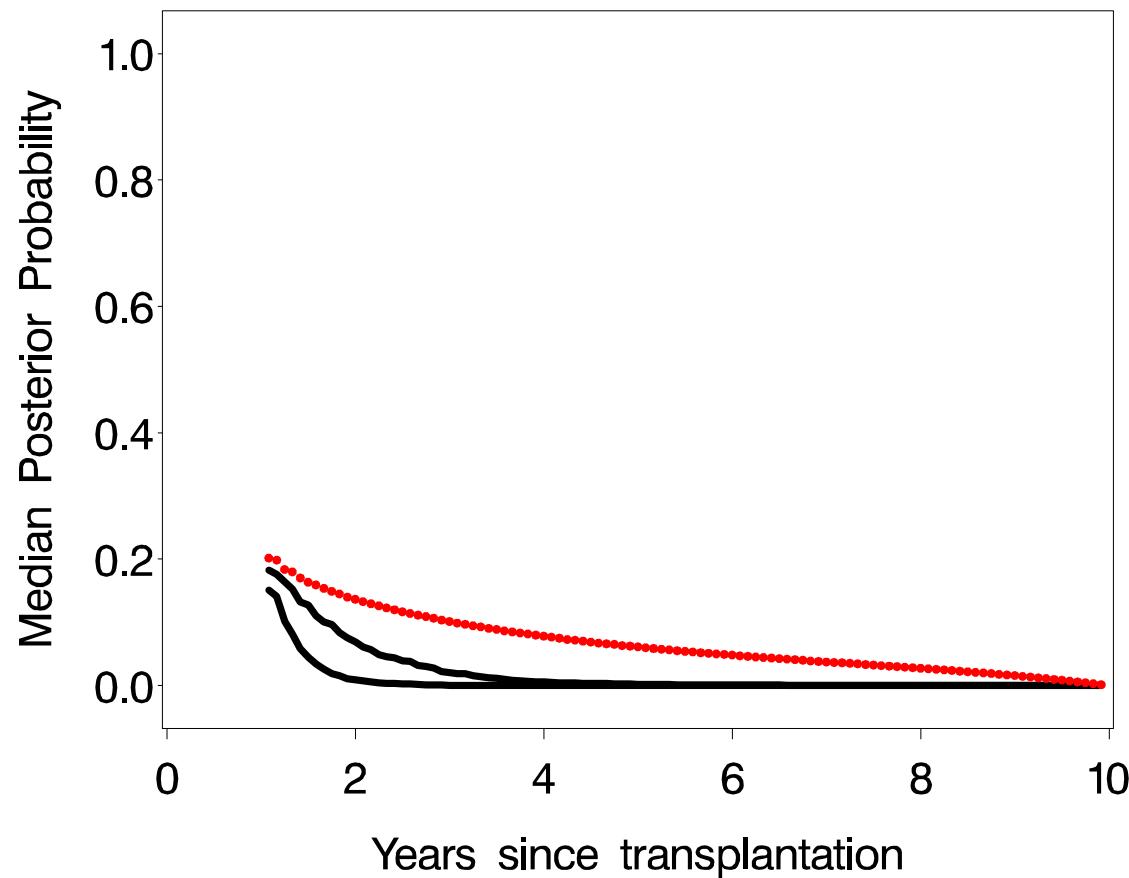
Median posterior probability: Non-failures

- Only GFR:



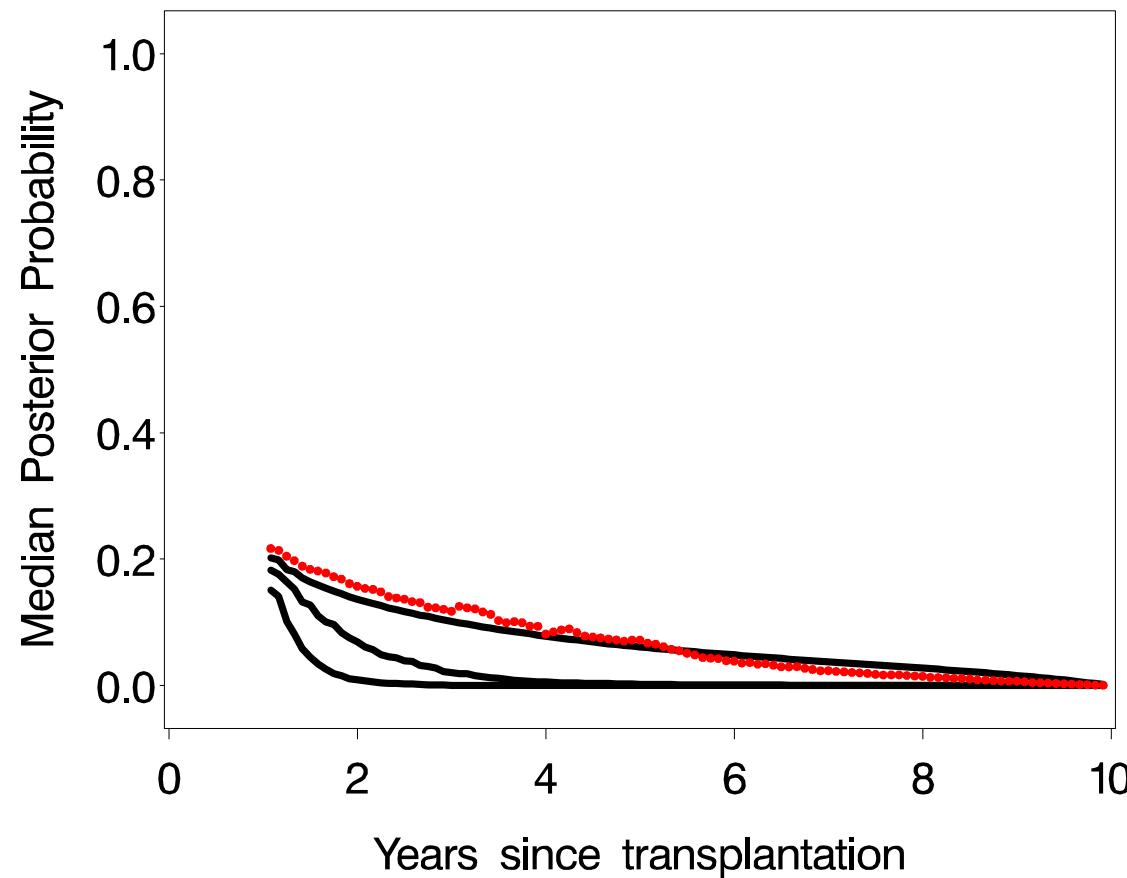
Median posterior probability: Non-failures

- Only Proteinuria:



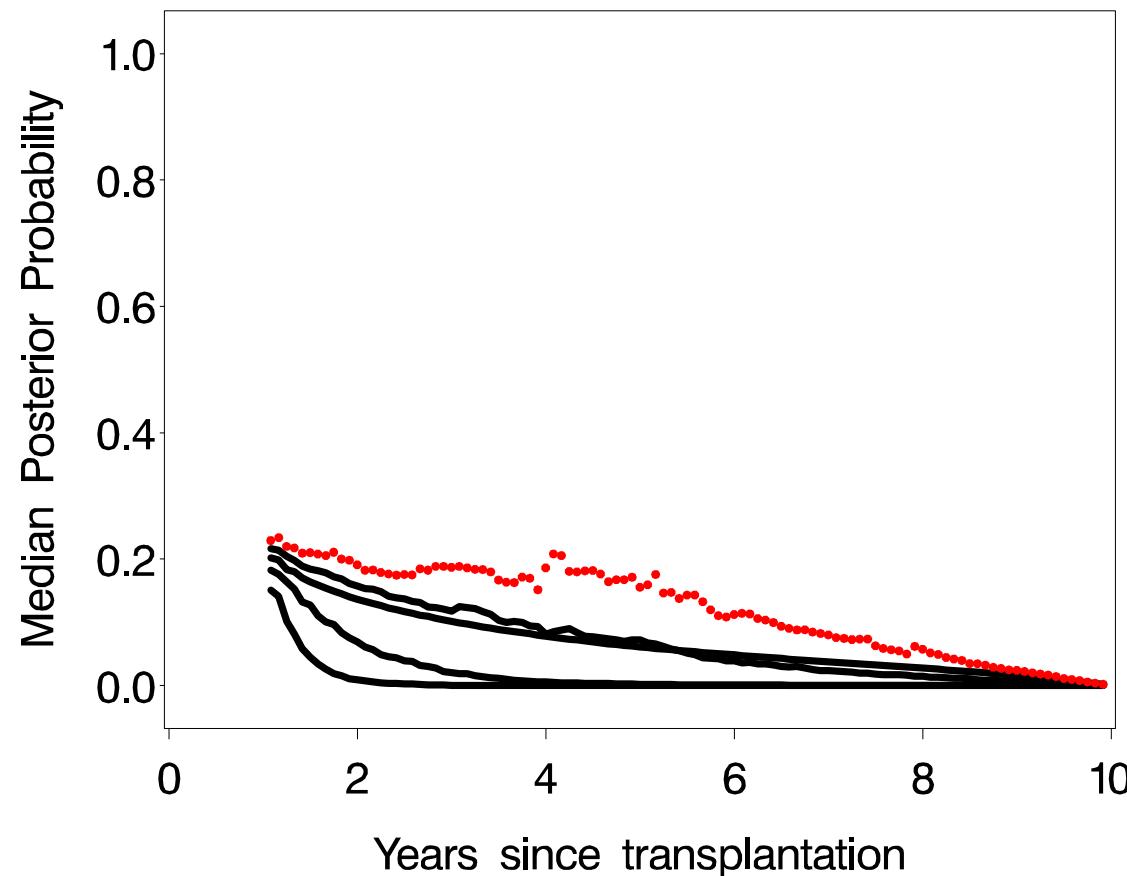
Median posterior probability: Non-failures

- Only Mean Blood Pressure:



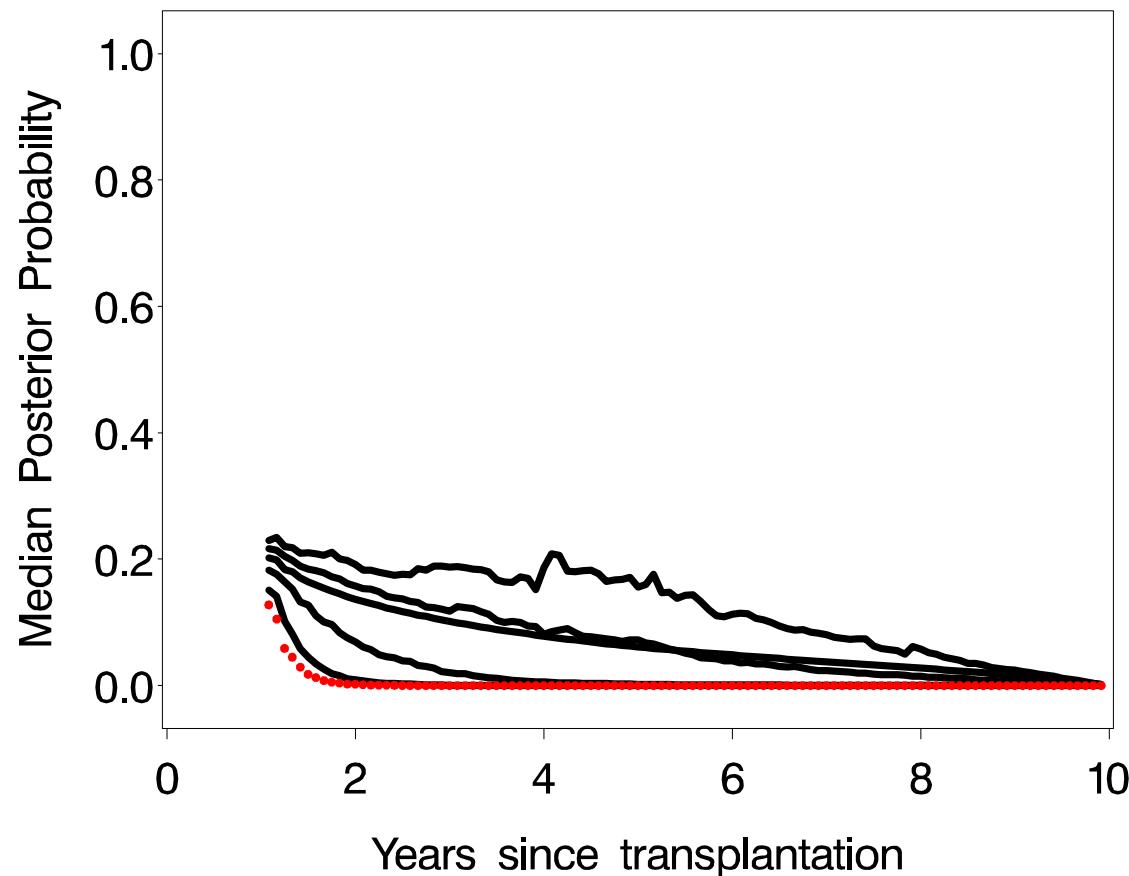
Median posterior probability: Non-failures

- Highest posterior probability over all four markers:



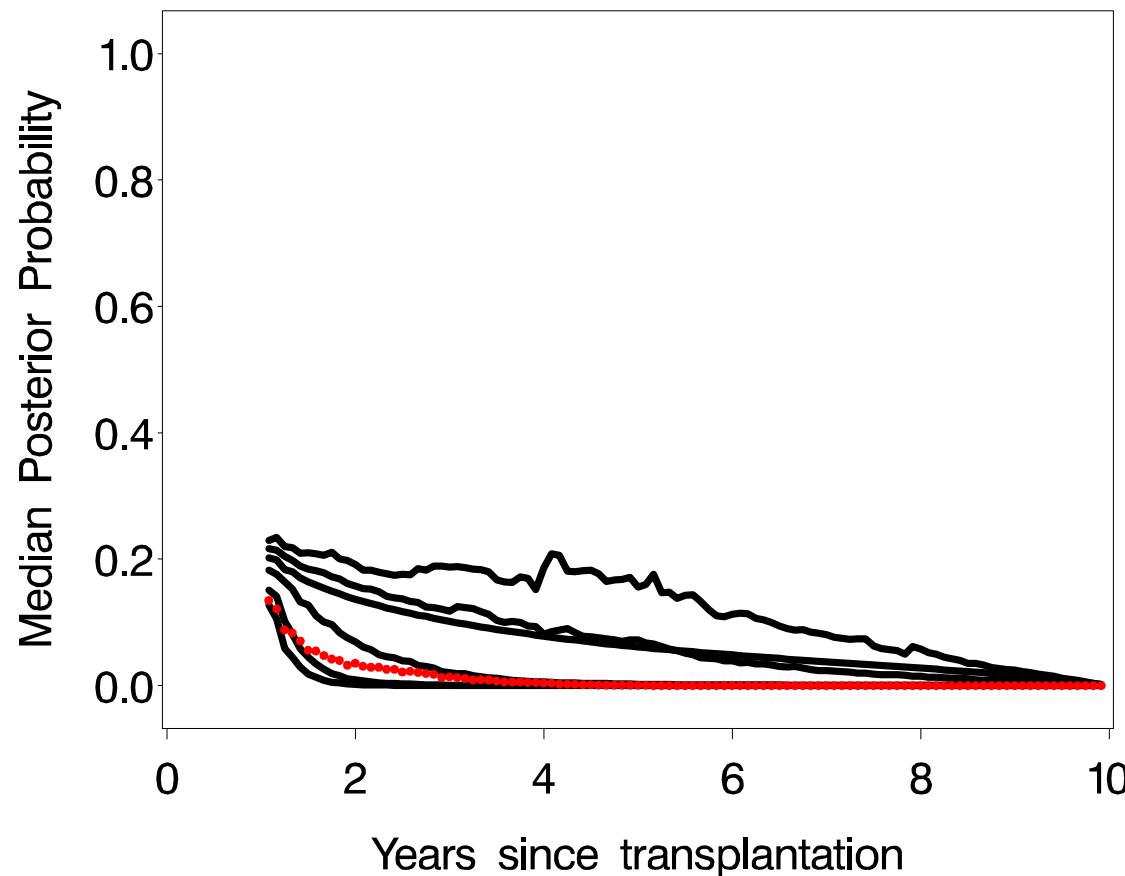
Median posterior probability: Non-failures

- All four markers, using joint model with uncorrelated markers:



Median posterior probability: Non-failures

- All four markers, using joint model with **correlated** markers:



Conclusions

- Discriminant analysis based on many outcomes, measured longitudinally, in an unbalanced design, is technically possible
- Allowing the longitudinal markers to be correlated considerably improves predictions
- A pattern-mixture approach allows for continuous updating of posterior probabilities
- Various mixed models can be combined and fitted using pairwise fitting approach



The End !